

Customer Profile Insights

描述性分析

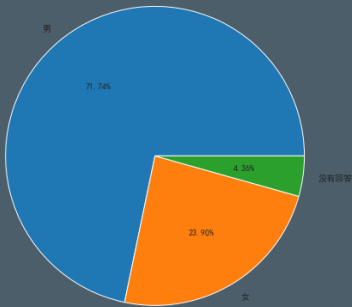
#### 数据缺失值百分比 ####	车主id	3955	non-null	object	
车主id	0.000000	姓名	3955	non-null	object
姓名	0.000000	性别	3947	non-null	object
性别	0.000000	生日	1965	non-null	datetime64[ns]
生日	47.599223	教育程度	2306	non-null	object
教育程度	41.659728	职位	2473	non-null	object
职位	37.246739	行业	2152	non-null	object
行业	45.628643	家庭年收入	1588	non-null	object
家庭年收入	59.089648	家庭成员人数	1138	non-null	object
家庭成员人数	72.356370	是否拥有驾照	1059	non-null	object
是否拥有驾照	74.687760	是否大客户	30	non-null	object
手机	0.000000	手机	3954	non-null	float64
省份	0.000000	省份	3693	non-null	object
城市	0.000000	城市	3611	non-null	object
		址区	3086	non-null	object
dtype: float64		邮编	2647	non-null	object

客户基本数据有3955条，16个特征维度，其中7个特征严重缺失，经过缺失值填补和特征转换，最终的到3602条有效数据，以及筛选出10条相关性较高的特征。

性别	教育程度	职位	行业	家庭年收入	家庭成员人数	是否拥有驾照	省份	城市	年龄	
车主ID										
1-8MN-68814	男	初中	总裁/总经理/总监/企业高管	家居、装饰	500000.0	4.0	Y	北京市	北京市	43.695031
1-3YPW-6314	男	本科	公司拥有者（老板）/合伙人	家居、装饰	500000.0	4.0	Y	江西省	景德镇市	40.000000
1-8LJ-21917	女	高中	总裁/总经理/总监/企业高管	电气、电器、仪器制造行业	100000.0	3.0	Y	河北省	衡水市	51.333333
1-8MK-70795	女	高中	总裁/总经理/总监/企业高管	电气、电器、仪器制造行业	230000.0	3.0	Y	安徽省	合肥市	57.000000
1-3S29-22819	男	高中	总裁/总经理/总监/企业高管	电气、电器、仪器制造行业	230000.0	3.0	Y	山东省	菏泽市	37.000000

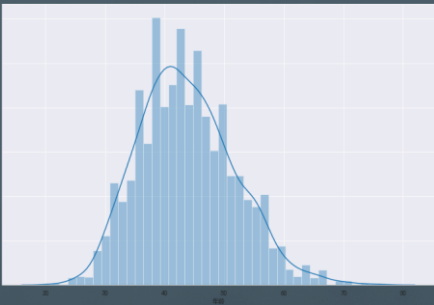
特征分析

性别



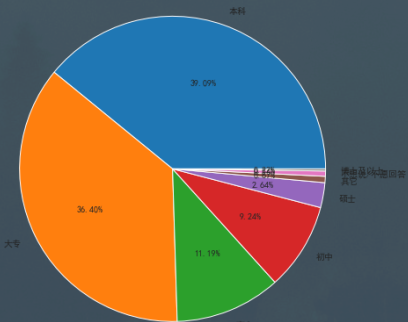
客户主要集中在男性群体，占比达到75.7%；说明我们更应注重男性偏好的相关车型的业务发展。

年龄



客户年龄分布集中在35到50岁，峰值在40岁左右；说明客户集中在中年人群，更适合推荐中年人偏好的车型。

教育程度



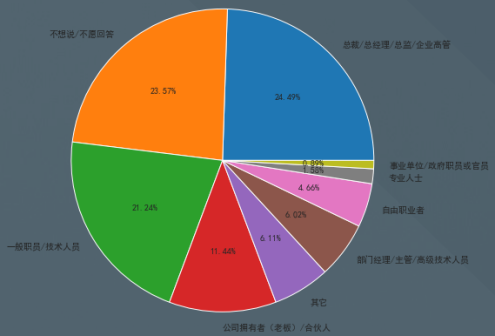
38.11%的客户是本科学历，36.68%的客户是大专学历，11.42%的客户是高中学历。硕士、博士级以上只占了3%左右；说明客户的教育程度在中等水平。

Customer Profile Insights

特征分析

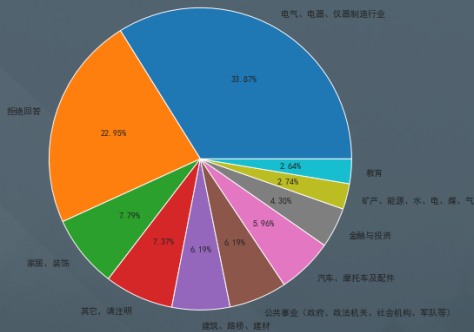
职位

24.5%左右的客户是企业高管等高层职位，有21.23%的客户是一般职员或者技术人员，而事业单位或者正负官员占比是最少的，不到1%；并且客户的职位呈现两级分化，高层职位和低层职位占比相近。



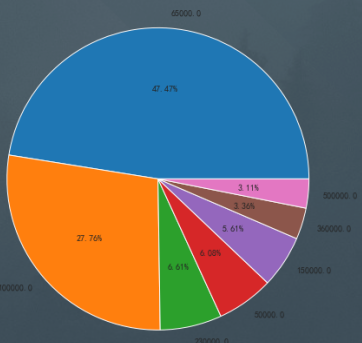
行业

33.85%的客户都集中在电气、电源、仪器制造行业。而其他行业。例如：家居、建筑、汽车、政府、金融、教育等，比重都差距不大，其中也有23%的人拒绝回答自己从事的行业。



教育程度

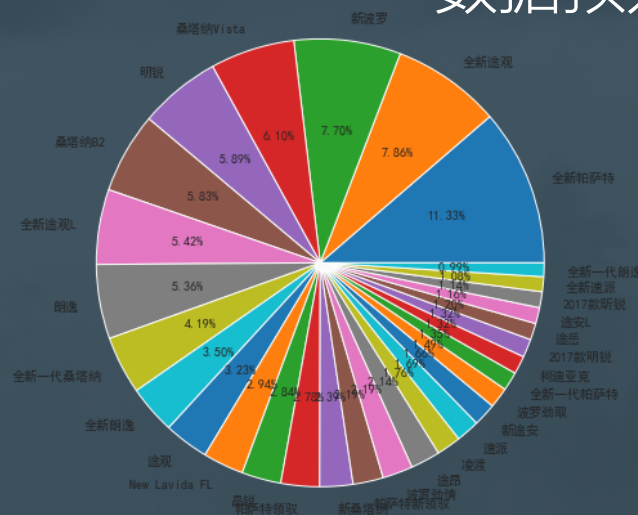
47.5%的客户收入水平在65000元左右，27.7%在8-12万元；年收入5万元以下和年收入在12-18万或18-27万的客户基本持平，都在6%左右；客户普遍集中在中等偏低收入水平的客户，总共占比大概在86%左右（针对5-15万年收入群体）。



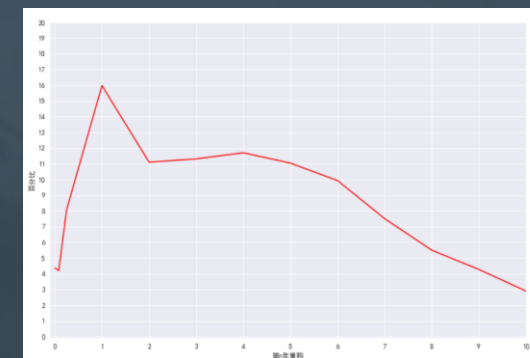
Repurchase Related Insights

数据预处理

```
#### 数据缺失值百分比 ####
车主ID          0.000000
VVIN            2.231136
市场车型        0.501379
车型            0.501379
购买日期        0.250689
开票价格        0.000000
ASSET_REF_EXPR  0.501379
OU_NAME         6.442717
OU_ABREV        6.442717
OU_CITY         6.492855
OU_COUNTY       9.275508
```



价格差异	
count	3.343000e+03
mean	5.846422e+04
std	1.146036e+05
min	-4.815000e+05
25%	0.000000e+00
50%	4.980000e+04
75%	1.265000e+05
max	2.222900e+06



数据概况&缺失值处理

客户重购数据含有7978行和11列，缺失值占比均低于10%，并且有486位客户暂时还未重购，而3343位客户有重购记录。

一共有76中不同款车型，平均价格在183000元上下。

市场车型分析

占据主导地位的车款是全新帕萨特（11.33%）；最贵的车款是途锐，均价在639000元左右，我根据对价格的分析结果对其进行评级（共7个等级），所以针对不同类型的客户，我们根据其基本信息、购买力、消费意愿以及积极程度来预测他们的回购情况，并相应的推荐车款。

重购价差分析

总体来说，重购车辆的价格会比首购平均高出46000元左右，变化量的中位数在7000元左右，最大价格差在220左右，最小则为-48万左右，差异变动很大，但整体成正态分布。

重购时间差分析

上图为重购情况和重购时间差的曲线图；当天重购和当月重购的比率接近（4%），而大部分客户选择在第一年内重购（16%），在第二、三、四、五年里重购的比重也十分相似（11%），之后重购率稳步下降；所以客户集中在短中期（5年以内）回购。

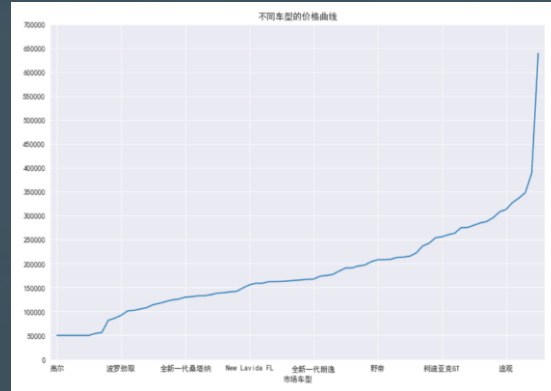
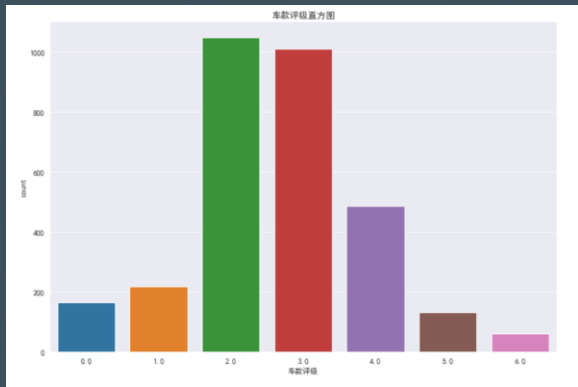
Repurchase Related Insights

数据分析&挖掘

	开票价格	价格差异	时间差异	车款评级
车主ID				
1-10LD-2240	256600.0	-8400.0	6.0	4.0
1-10ZA-842	131900.0	-163800.0	510.0	2.0
1-11GN3TJ	119650.0	-1700.0	495.0	2.0
1-11GNTKY	191850.0	-143900.0	3.0	3.0
1-11HMNQQ	323850.0	63900.0	493.0	5.0
...
1-Z7U-601	187950.0	24100.0	947.0	3.0
1-ZG3W1L	224900.0	76000.0	716.0	3.0
1-ZPOITB	114900.0	4000.0	39.0	2.0
1-ZUL-1836	192250.0	95100.0	1403.0	3.0
1-ZUZ-86	197880.0	26560.0	172.0	3.0

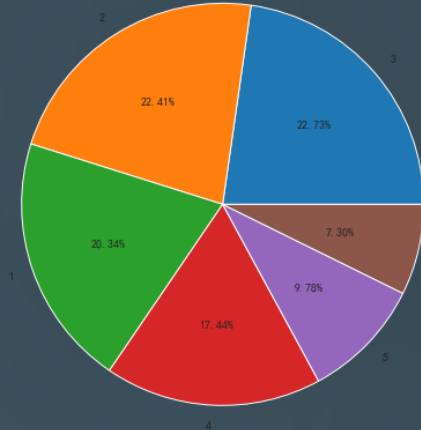
回购客户特征分析

开票价格、回购价格差异、回购时间差异对于车款评级具有决定性因素，车款评级分布如下：



客户活跃度分析

根据客户的消费能力和重购时差来界定客户的活跃程度；65%的客户活跃度较高，下图为分级情况。



市场车型	3120 non-null	object
开票价格	3120 non-null	float64
品牌	3120 non-null	object
所购买省份	3120 non-null	object
所购买公司	3120 non-null	object
价格差异	3120 non-null	float64
时间差异	3120 non-null	float64
车款评级	3120 non-null	float64
购买年份	3120 non-null	int64
购买月份	3120 non-null	int64
购买日	3120 non-null	int64
性别	3120 non-null	object
教育程度	3120 non-null	object
职位	3120 non-null	object
行业	3120 non-null	object
家庭年收入	3120 non-null	float64
家庭成员人数	3120 non-null	float64
是否拥有驾照	3120 non-null	object
省份	3120 non-null	object
城市	3120 non-null	object
年龄	3120 non-null	float64

客户信息交互

通过连接客户基本信息数据和客户回购信息数据进行信息交互，整合出了3120条有效数据以及21个有效特征；数据概况入下：

	市场车型	开票价格	品牌	所购省份	所购公司	价格差异	时间差异	车款评级	购买年份	购买月份	教育程度	性别	职位	行业	家庭年收入	家庭成员人数	是否拥有驾照	省份	城市	年龄
车主ID																				
1-10LD-2240	全新帕萨特	256600.0	大众	河北省	沧州兴顺汽车销售服务有限公司	-8400.0	6.0	4.0	2013	12	16	男	高中	公司所有者(老板) / 合伙人	汽车、摩托车及配件	100000.0	3.0	Y	河北省	38.00
1-10ZA-842	全新朗逸	131900.0	大众	江苏省	南京公用发展股份有限公司	-163800.0	510.0	2.0	2015	5	6	女	大学	公司所有者(老板) / 合伙人	电气、电器、仪器制造业	230000.0	3.0	Y	江苏省	55.00
1-11GN3TJ	全新桑塔纳	119650.0	大众	安徽省	滁州德胜汽车销售服务有限公司	-1700.0	495.0	2.0	2017	7	8	男	高中	公司所有者(老板) / 合伙人	其它、制造业	500000.0	2.0	Y	安徽省	60.33
1-11GNTKY	全新一代帕萨特	191850.0	大众	北京市	北京市艾迪汽车销售有限公司	-143900.0	3.0	3.0	2016	3	3	女	本科	一般职员/技术人员	电气、电器、仪器制造业	65000.0	3.0	Y	河北省	41.50
1-11HMNQQ	全新途观L	323850.0	大众	贵州省	贵州华通汽车销售服务有限公司	63900.0	493.0	5.0	2017	7	7	男	初中	公司所有者(老板) / 合伙人	建筑、建材、银行	65000.0	4.0	Y	贵州省	55.67

市场车型	开票价格	品牌	所购买省份	所购买公司	价格差异	时间差异	车款评级	购买年份	购买月份	购买省份	性别	职业	职位	行业	家庭年收入	家庭总资产人数	是否拥有驾照	省份	城市	年龄	
车主ID																					
1-18T5X0R	2,636,144	157250.0	0	2,808,284	2,630,323	10700.0	77.0	2.0	2016	11	14	1.0	5	3	2,688,165	65000.0	3.0	1	2,648,530	2,644,588	33.00
1-XOL588	3,188,376	266250.0	1	2,690,992	2,647,845	9300.0	200.0	4.0	2016	8	19	1.0	5	0	2,477,465	65000.0	3.0	1	2,636,150	2,582,867	33.00
1-40Y-51680	3,230,069	153150.0	1	2,575,592	2,674,970	206300.0	2799.0	2.0	2017	11	21	0.0	4	0	2,445,586	65000.0	3.0	1	2,674,882	2,625,331	49.00
1-3FZ-4739	2,692,676	128450.0	0	2,629,005	2,655,653	12900.0	102.0	2.0	2015	5	15	1.0	5	4	2,647,529	100000.0	3.0	1	2,639,664	2,636,107	35.33
1-80-1136	2,619,543	182150.0	1	2,628,727	2,649,557	-17300.0	1066.0	2.0	2016	4	6	1.0	5	0	2,464,505	65000.0	3.0	1	2,648,655	2,595,000	36.00

特征工程&相关性分析

对categorical变量进行encoding处理，对其进行适当的编码转换，从而为带入模型做准备；我采取了target encoding以及平滑处理解决了数据泄露以及维度爆炸的问题，并使用label encoding对序列数据进行编码，从而使得所有数据序列化、数字化。

下图为特征之间的关于heatmap的相关性分析图表。



